

CHAPTER 1

Foundations of Data Analysis and Python

Introduction

In today's data-rich landscape, data is much more than a collection of numbers or facts, it's a powerful resource that can influence decision-making, policy formation, product development, and scientific discovery. To turn these raw inputs into meaningful insights, we rely on statistics, the discipline dedicated to collecting, organizing, summarizing, and interpreting data. Statistics not only helps us understand patterns and relationships but also guides us in making evidence-based decisions with confidence. This chapter examines fundamental concepts at the heart of data analysis. We'll explore what data is and why it matters, distinguish between various types of data and their levels of measurement, and consider how data can be categorized as univariate, bivariate, or multivariate. We'll also highlight different data sources, clarify the roles of populations and samples, and introduce crucial data preparation tasks including cleaning, wrangling, and manipulation to ensure data quality and integrity.

For example, consider you have records of customer purchases at an online store everything from product categories and prices to transaction dates and customer demographics. Applying statistical principles and effective data preparation techniques to this information can reveal purchasing patterns, highlight which product lines drive the most revenue, and suggest targeted promotions that improve the shopping experience.

Structure

In this chapter, we will discuss the following topics:

- Environment setup
- Software installation
- Basic overview of technology
- Statistics, data, and its importance
- Types of data
- Levels of measurement
- Univariate, bivariate, and multivariate data
- Data sources, methods, population, and samples
- Data preparation tasks
- Wrangling and manipulation

Objectives

By the end of this chapter, readers will learn the basics of statistics and data, such as, what they are, why they are important, how they vary in type and application, and the basic data collection and manipulation techniques. Moreover, this chapter explains different level of measurements, data analysis techniques, its source, collection methods, their quality and cleaning. You will also learn how to work with data using Python, a powerful and popular programming language that offers many tools and

libraries for data analysis.

Environment setup

To set up the environment and to run the sample code for statistics and data analysis in Python, the three options are as follows:

- Download and install Python from <https://www.python.org/downloads/>. Other packages need to be installed explicitly on top of Python. Then, use any **integrated development environment (IDE)** like visual studio code to execute Python code.
- You can also use Anaconda, a Python distribution designed for large-scale data processing, predictive analytics, and scientific computing. The Anaconda distribution is the easiest way to code in Python. It works on Linux, Windows, and Mac OS X. It can be downloaded from <https://www.anaconda.com/distribution/>.
- You can also use cloud services, which is the easiest of all options but requires internet connectivity to use. Cloud providers like Microsoft Azure Notebooks, GitHub Code Spaces and Google Collaboratory are very popular. Following are a few links:
 - **Microsoft Azure Notebooks:** <https://notebooks.azure.com/>
 - **GitHub Codespaces:** Create a GitHub account from <https://github.com/join> then, once logged in, create a repository from <https://github.com/new>. Once the repository is created, open the repository in the codespace by using the following instructions:
<https://docs.github.com/en/codespaces/developing-in-codespaces/creating-a-codespace-for-a->

repository.

- **Google Collaboratory:** Create a Google account, open <https://colab.research.google.com/>, and create a new notebook.

Azure Notebook GitHub Codespace and Google Collaboratory are cloud-based and easy-to-use platforms. To run and set up an environment locally, install the Anaconda distribution on your machine and follow the software installation instructions.

Software installation

Now, let us look at the steps to install Anaconda to run the sample code and tutorials on the local machine as follows:

1. Download the Anaconda Python distribution from the following link: <https://www.anaconda.com/download>
2. Once the download is complete, run the setup to begin the installation process.
3. Once the Anaconda application has been installed, click **Close** and move to the next step to launch the application.

Check Anaconda installation instructions in the following:

<https://docs.anaconda.com/free/anaconda/install/index.html>

Launch application

Now, let us launch the installed Anaconda navigator and the JupyterLab in it.

Following are the steps:

1. After installing the Anaconda navigator, open any Anaconda navigator and then install and launch JupyterLab.
2. This will start the Jupyter server listening on port 8888.

Usually, a pop-up window comes with a default browser, but you can also start the JupyterLab application on any web browser, Google Chrome preferred, and go to the following URL:

<http://localhost:8888/>

3. A blank notebook is launched in a new window. You can write Python code on it.

4. Select the cell and press run to execute the code.

The environment is now ready to write, run and execute tutorials.

Basic overview of technology

Python, NumPy, pandas, Sklearn, Matplotlib will be used in most of the tutorials. Let us have a look at them in the following section.

Python

To know more about Python and installation you can refer to the following link:

<https://www.python.org/about/gettingstarted/>. Execute Python-version in terminal or command prompt, and if you see the Python version as output, you are good to go, else, install Python. There are different ways to install Python packages on Jupyter Notebook, depending on the package manager you use and the environment you work in, as follows:

- If you use **pip** as your package manager, you can install packages directly from a code cell in your notebook by typing **!pip install <package_name>** and running the cell. Then replace **<package_name>** with the name of the package you want to install.
- If you use **conda** as your package manager, you can

install packages from a JupyterLab cell by typing **!conda install <package_name> --yes** and running the cell. The **--yes** flag is to avoid prompts that asks for confirmation.

- If you want to install a specific version of Python for your notebook, you can use the ipykernel module to create a new kernel with that version. For example, if you have *Python 3.11* and pip installed on your machine, you can type **!pip3.11 install ipykernel** and **!python3.11 -m ipykernel install -user** in two separate code cells and run them. Then, you can select *Python 3.11* as your kernel from the kernel menu.

Further tutorials will be based on the JupyterLab.

pandas

pandas is mainly used for data analysis and manipulation in Python. More can be read at:

<https://pandas.pydata.org/docs/>

Following are the ways to install pandas:

- In Jupyter Notebook, execute **pip install pandas**
- In the conda environment, execute **conda install pandas --yes**

NumPy

NumPy is a Python package for numerical computing, multi-dimensional array, and math computation. More can be read at <https://pandas.pydata.org/docs/>.

Following are the ways to install NumPy:

- In Jupyter Notebook, execute **pip install pandas**
- In the conda environment, execute **conda install pandas -yes**

Sklearn

Sklearn is a Python package that provides tools for machine learning, such as data preprocessing, model selection, classification, regression, clustering, and dimensionality reduction. Sklearn is mainly used for predictive data analysis and building machine learning models. More can be read at <https://scikit-learn.org/0.21/documentation.html>.

Following are the ways to install Sklearn:

- In Jupyter Notebook, execute **pip install scikit-learn**
- In the conda environment, execute **conda install scikit-learn -yes**

Matplotlib

Matplotlib is mainly used to create static, animated, and interactive visualizations (plots, figures, and customized visual style and layout) in Python. More can be read at <https://matplotlib.org/stable/index.html>.

Following are the ways to install Matplotlib:

- In Jupyter notebook, excute **pip install matplotlib**
- In the conda environment, execute **conda install matplotlib --yes**

Statistics, data and its importance

As established, statistics is a disciplined approach that enables us to derive insights from diverse forms of data. By applying statistical principles, we can understand what is happening around us quantitatively, evaluate claims, avoid misleading interpretations, produce trustworthy results, and support data-driven decision-making. Statistics also equips us to make predictions and deepen our understanding of the subjects we study.

Data, in turn, serves as the raw material that fuels statistical analysis. It may take various forms—numbers, words,

images, sounds, or videos—and provides the foundational information needed to extract useful knowledge and generate actionable insights. Through careful examination and interpretation, data leads us toward new discoveries, informed decisions, and credible forecasts.

Ultimately, data and statistics are interdependent. Without data, statistics has no basis for drawing conclusions; without statistics, raw data remains untapped and lacks meaning. When combined, they answer fundamental WH questions—Who, What, When, Where, Why, and How—with clarity and confidence, guiding our understanding and shaping the decisions we make.

Types of data

Data can be in different form and type but generally it can be divided into two types, that is, qualitative and quantitative.

Qualitative data

Qualitative data cannot be measured or counted in numbers. Also known as categorical data, it is descriptive, interpretation-based, subjective, and unstructured. It describes the qualities or characteristics of something. It helps to understand the reasoning behind it by asking why, how, or what. It includes nominal and ordinal data. For example, gender of person, race of a person, smartphone brand, hair color type, marital status, and occupation of a person.

Tutorial 1.1: To implement creating a data frame consisting of only qualitative data.

To create a data frame with pandas, **import pandas as pd**, then use the **DataFrame()** function and pass a data source, such as a dictionary, list, or array, as an argument.


```

1. #Import the pandas library to create a pandas Dataframe
2. import pandas as pd
3. # Sample qualitative data
4. qualitative_data = {
5.     'Name': ['John', 'Alice', 'Bob', 'Eve', 'Michael'],
6.     'City': ['New York', 'Los Angeles', 'Chicago', 'San Francisco', 'Miami'],
7.     'Gender': ['Male', 'Female', 'Male', 'Female', 'Male'],
8.     'Occupation': ['Engineer', 'Artist', 'Teacher', 'Doctor', 'Lawyer'],
9.     'Race': ['Black', 'White', 'Asian', 'Indian', 'Mongolian'],
10.    'Smartphone Brand': ['Apple', 'Samsung', 'Xiaomi', 'Apple', 'Google']
11. }
12. # Create the DataFrame
13. qualitative_df = pd.DataFrame(qualitative_data)
14. # Prints the created DataFrame
15. print(qualitative_df)

```

Output:

```

1.   Name  City  Gender Occupation  Race
   Smartphone Brand
2. 0  John  New York  Male   Engineer  Black
   Apple
3. 1  Alice  Los Angeles  Female  Artist   White
   Samsung
4. 2  Bob   Chicago   Male   Teacher  Asian   Xiaomi
5. 3  Eve   San Francisco  Female  Doctor   Indian
   Apple
6. 4  Michael  Miami   Male   Lawyer
   Mongolian  Google

```

Row consisting of numbers 0, 1, 2, 3, and 4 is the index column, not part of the qualitative data. To exclude it from output, hide the index column using **to_string()** as follows:

```
1. print(qualtitative_df.to_string(index=False))
```

Output:

1.	Name	City	Gender	Occupation	Race	Smartphone Brand
2.	John	New York	Male	Engineer	Black	Apple
3.	Alice	Los Angeles	Female	Artist	White	Samsung
4.	Bob	Chicago	Male	Teacher	Asian	Xiaomi
5.	Eve	San Francisco	Female	Doctor	Indian	Apple
6.	Michael	Miami	Male	Lawyer	Mongolian	Google

While we often think of data in terms of numbers, many other forms such as images, audio, videos, and text they can also represent quantitative information when suitably encoded (e.g., pixel intensity values in images, audio waveforms, or textual features like word counts).

Tutorial 1.2: To implement accessing and creating a data frame consisting of the image data.

In this tutorial, we'll work with the open-source Olivetti faces dataset, which consists of grayscale face images collected at AT&T Laboratories Cambridge between April 1992 and April 1994. Each face is represented by numerical pixel values, making them a form of quantitative data. By organizing this data into a DataFrame, we can easily manipulate, analyze, and visualize it for further insights.

To create a data frame consisting of the *Olivetti faces*

dataset, you can use the following steps:

1. Fetch the **Olivetti faces** dataset from sklearn using the **sklearn.datasets.fetch_olivetti_faces** function. This will return an object that holds the data and some metadata.
2. Use the **pandas.DataFrame** constructor to create a data frame from the data and the feature names. You can also add a column for the target labels using the target and **target_names** attributes of the object.
3. Use the pandas method to display and analyze the data frame. For example, you can use **df.head()**, **df.describe()**, **df.info()**.

```
1. import pandas as pd
2. # Import datasets from the sklearn library
3. from sklearn import datasets
4. # Fetch the Olivetti faces dataset
5. faces = datasets.fetch_olivetti_faces()
6. # Create a dataframe from the data and feature names
7. df = pd.DataFrame(faces.data)
8. # Add a column for the target labels
9. df["target"] = faces.target
10. # Display the first 3 rows of the dataframe
11. print(f"{df.head(3)}")
12. # Print new line
13. print("\n")
14. # Display the first image in the dataset
15. import matplotlib.pyplot as plt
16. plt.imshow(df.iloc[0, :-1].values.reshape(64, 64), cmap="gray")
17. plt.title(f"Image of person {df.iloc[0, -1]}")
```

18. plt.show()

Quantitative data

Quantitative data is measurable and can be expressed numerically. It is useful for statistical analysis and mathematical calculations. For example, if you inquire about the number of books people have read in a month, their responses constitute quantitative data. They may reveal that they have read, let us say, **three books, zero books, or ten books**, providing information about their reading habits. Quantitative data is easily comparable and allows for calculations. It can provide answers to questions such as **How many? How much? How often?** and **How fast?**

Tutorial 1.3: To implement creating a data frame consisting of only quantitative data is as follows:

```
1. #Import the pandas library to create a pandas DataFrame
   e
2. import pandas as pd
3. quantitative_df = pd.DataFrame({
4.     "price": [300000, 250000, 400000, 350000, 450000],
5.     "distance": [10, 15, 20, 25, 30],
6.     "height": [170, 180, 190, 160, 175],
7.     "weight": [70, 80, 90, 60, 75],
8.     "salary": [5000, 6000, 7000, 8000, 9000],
9.     "temperature": [25, 30, 35, 40, 45],
10. })
11. # Print the DataFrame without index
12. print(quantitative_df.to_string(index=False))
```

Output:

```
1. price distance height weight salary temperature
2. 300000      10    170     70    5000          25
```

3.	250000	15	180	80	6000	30
4.	400000	20	190	90	7000	35
5.	350000	25	160	60	8000	40
6.	450000	30	175	75	9000	45

Tutorial 1.4: To implement accessing and creating a data frame by loading the tabular **iris** data.

Iris tabular dataset contains 150 samples of iris flowers with four features, that is, sepal length, sepal width, petal length, and petal width and three classes, that is, setosa, versicolor, and virginica. The sepal length, sepal width, petal length, petal width, and target (class) are columns of the table¹.

To create a data frame consisting of the **iris** dataset, you can use the following steps:

1. First, you need to load the **iris** dataset from sklearn using the **sklearn.datasets.load_iris** function. This will return a bunch object that holds the data and some metadata.
2. Next, you can use the **pandas.DataFrame** constructor to create a data frame from the data and the feature names. You can also add a column for the target labels using the target and **target_names** attributes of the bunch object.
1. Finally, you can use the panda method to display and analyze the data frame. For example, you can use **df.head()**, **df.describe()**, **df.info()** as follows:

```

1. import pandas as pd
2. # Import dataset from sklearn
3. from sklearn import datasets
4. # Load the iris dataset
5. iris = datasets.load_iris()
6. # Create a dataframe from the data and feature nam

```

```

es
7. df = pd.DataFrame(iris.data, columns=iris.feature_names)
8. # Add a column for the target labels
9. df["target"] = iris.target
10. # Display the first 5 rows of the dataframe
11. df.head()

```

Level of measurement

Level of measurement is a way of classifying data based on how precise it is and what we can do with it. Generally, they are four, that is, nominal, ordinal, interval and ratio. Nominal is a category with no inherent order, such as colors. Ordinal is a category with a meaningful order, such as education levels. Interval is equal intervals but no true zero, such as temperature in degrees Celsius, and ratio are equal intervals with a true zero, such as age in years.

Nominal data

Nominal data is qualitative data that does not have a natural ordering or ranking. For example, gender, religion, ethnicity, color, brand ownership of electronic appliances, and person's favorite meal.

Tutorial 1.5: To implement creating a data frame consisting of qualitative nominal data, is as follows:

```

1. #Import the pandas library to create a pandas Dataframe
2. import pandas as pd
3. nominal_data = {
4.     "Gender": ["Male", "Female", "Male", "Female", "Male"],
5.     "Religion": ["Hindu", "Muslim", "Christian", "Buddhist"]

```

```

    ", "Jewish"],
6.    "Ethnicity": ["Indian", "Pakistani", "American", "Chinese", "Israeli"],
7.    "Color": ["Red", "Green", "Blue", "Yellow", "White"],
8.    "Electronic Appliances Ownership": ["Samsung", "LG", "Apple", "Huawei", "Sony"],
9.    "Person Favorite Meal": ["Biryani", "Kebab", "Pizza", "Noodles", "Falafel"],
10.   "Pet Preference": ["Dog", "Cat", "Parrot", "Fish", "Hamster"]
11. }
12. # Create the DataFrame
13. nominal_df = pd.DataFrame(nominal_data)
14. # Display the DataFrame
15. print(nominal_df)

```

Output:

```

1. Gender  Religion  Ethnicity  Color  Electronic Appliance
   s Ownership \
2. 0  Male    Hindu    Indian    Red                      Samsu
   ng
3. 1  Female  Muslim  Pakistani  Grezn
   LG
4. 2  Male  Christian  American  Blue                      App
   le
5. 3  Female  Buddhist  Chinese  Yellow                      Hu
   awai
6. 4  Male    Jewish  Israeli  Whie                      Sony

7.
8. Person Favorite Meal Pet Preference
9. 0                      Biryani          Dog
10. 1                      Kebab           Cat

```

11. 2	Pizza	Parrot
12. 3	Noodles	Fish
13. 4	Falafel	Hamster

Ordinal data

Ordinal data is qualitative data that has a natural ordering or ranking. For example, student ranking in class (1st, 2nd, or 3rd), educational qualification (high school, undergraduate, or graduate), satisfaction level (bad, average, or good), income level range, level of agreement (agree, neutral, or disagree).

Tutorial 1.6: To implement creating a data frame consisting of qualitative ordinal data is as follows:

```

1. import pandas as pd
2. ordinal_data = {
3.     "Student Rank in a Class": ["1st", "2nd", "3rd", "4th",
4.     "5th"],
5.     "Educational Qualification": ["Graduate", "Undergrad
6.     uate", "High School", "Graduate", "Undergraduate"],
7.     "Satisfaction Level": ["Good", "Average", "Bad", "Aver
8.     age", "Good"],
9.     "Income Level Range": ["80,000-100,000", "60,000-
10.    80,000", "40,000-60,000", "100,000-120,000", "50,000-
11.    70,000"],
12.     "Level of Agreement": ["Agree", "Neutral", "Disagree"
13.     , "Neutral", "Agree"]
14. }
15. ordinal_df = pd.DataFrame(ordinal_data)
16. print(ordinal_df)

```

Output:

```

1. Student Rank in a Class Educational Qualification Sati
2. sfaction Level \
3. 0 1st Graduate Good

```


3.	1	2nd	Undergraduate	Average
4.	2	3rd	High School	Bad
5.	3	4th	Graduate	Average
6.	4	5th	Undergraduate	Good
7.				
8.	Income Level Range Level of Agreement			
9.	0	80,000-100,000	Agree	
10.	1	60,000-80,000	Neutral	
11.	2	40,000-60,000	Disagree	
12.	3	100,000-120,000	Neutral	
13.	4	50,000-70,000	Agree	

Discrete data

Discrete data is quantitative data, integers or whole numbers, they cannot be subdivided into parts. For example, total number of students present in a class, cost of a cell phone, number of employees in a company, total number of players who participated in a competition, days in a week, number of books in a library, etc. For example, number of coins in a jar, it can only be a whole number like 1,2,3 and so on.

Tutorial 1.7: To implement creating a data frame consisting of quantitative discrete data is as follows:

```

1. import pandas as pd
2. discrete_data = {
3.     "Students": [25, 30, 35, 40, 45],
4.     "Cost": [500, 600, 700, 800, 900],
5.     "Employees": [100, 150, 200, 250, 300],
6.     "Players": [50, 40, 30, 20, 10],
7.     "Week": [7, 7, 7, 7, 7]
8. }
```

```
9. discrete_df = pd.DataFrame(discrete_data)
10. discrete_df
```

Output:

```
1. Students Cost Employees Players Week
2. 0 25 500 100 50 7
3. 1 30 600 150 40 7
4. 2 35 700 200 30 7
5. 3 40 800 250 20 7
6. 4 45 900 300 10 7
```

Continuous data

Continuous data is quantitative data that can take any value (including fractional value) within a range and have no gaps between them. No gaps mean that if a person's height is 1.75 meters, there is always a possibility of height being between 1.75 and 1.76 meters, such as 1.751 or 1.755 meters.

Interval data

Interval data is quantitative numerical data with inherent order. They always have an arbitrary zero, an arbitrary zero meaning no meaningful zero, chosen by convention, not by nature. For example, a temperature of zero degrees Fahrenheit does not mean that there is no heat or temperature, here, zero is an arbitrary zero point. For example, temperature (Celsius or Fahrenheit), GMAT score (200-800), SAT score (400-1600).

Tutorial 1.8: To implement creating a data frame consisting of quantitative interval data is as follows:

```
1. import pandas as pd
2. interval_data = {
3.     "Temperature": [10, 15, 20, 25, 30],
4.     "GMAT_Score": [600, 650, 700, 750, 800],
```

```

5.     "SAT_Score (400 - 1600)": [1200, 1300, 1400, 1500, 1600],
6.     "Time": ["9:00", "10:00", "11:00", "12:00", "13:00"]
7. }
8. interval_df = pd.DataFrame(interval_data)
9. # Print DataFrame as it is without print() also
10. interval_df

```

Output:

```

1. Temperature  GMAT_Score  SAT_Score (400 - 1600)  Time
2. 0  10  600  1200  9:00
3. 1  15  650  1300  10:00
4. 2  20  700  1400  11:00
5. 3  25  750  1500  12:00
6. 4  30  800  1600  13:00

```

Ratio data

Ratio data is naturally, numerical ordered data with an absolute, where zero is not arbitrary but meaningful. For example, height, weight, age, tax amount has true zero point that is fixed by nature, and they are measured on a ratio scale. Zero height means no height at all, like a point in space. There is nothing shorter than zero height. Zero tax amount means no tax at all, like being exempt. There is nothing lower than zero tax amount.

Tutorial 1.9: To implement creating a data frame consisting of quantitative ratio data is as follows:

```

1. import pandas as pd
2. ratio_data = {
3.     "Height": [170, 180, 190, 200, 210],
4.     "Weight": [60, 70, 80, 90, 100],
5.     "Age": [20, 25, 30, 35, 40],

```

```
6.     "Speed": [80, 90, 100, 110, 120],
7.     "Tax Amount": [1000, 1500, 2000, 2500, 3000]
8. }
9. ratio_df = pd.DataFrame(ratio_data)
10. ratio_df
```

Output:

```
1. Height Weight Age Speed Tax Amount
2. 0 170 60 20 80 1000
3. 1 180 70 25 90 1500
4. 2 190 80 30 100 2000
5. 3 200 90 35 110 2500
6. 4 210 100 40 120 3000
```

Tutorial 1.10: To implement loading the ratio data in a JSON format and displaying it.

Sometimes, data can be in **JSON**. The data used in the following Tutorial 1.10 is in **JSON** format. In that case **json.loads()** method can load it. **JSON** is a text format for data interchange based on JavaScript as follows:

```
1. # Import json
2. import json
3. # The JSON string:
4. json_data = """
5. [
6. {
7.     "Height": 170,
8.     "Weight": 60,
9.     "Age": 20,
10.    "Speed": 80,
11.    "Tax Amount": 1000
12. },
13. {
```

```
14.     "Height": 180,
15.     "Weight": 70,
16.     "Age": 25,
17.     "Speed": 90,
18.     "Tax Amount": 1500
19. },
20. {
21.     "Height": 190,
22.     "Weight": 80,
23.     "Age": 30,
24.     "Speed": 100,
25.     "Tax Amount": 2000
26. },
27. {
28.     "Height": 200,
29.     "Weight": 90,
30.     "Age": 35,
31.     "Speed": 110,
32.     "Tax Amount": 2500
33. },
34. {
35.     "Height": 210,
36.     "Weight": 100,
37.     "Age": 40,
38.     "Speed": 120,
39.     "Tax Amount": 3000
40. }
41. ]
42. ""
43. # Convert to Python object (list of dicts):
```

```
44. data = json.loads(json_data)
```

```
45. data
```

Output:

1. [{ 'Height': 170, 'Weight': 60, 'Age': 20, 'Speed': 80, 'Tax Amount': 1000},
2. { 'Height': 180, 'Weight': 70, 'Age': 25, 'Speed': 90, 'Tax Amount': 1500},
3. { 'Height': 190, 'Weight': 80, 'Age': 30, 'Speed': 100, 'Tax Amount': 2000},
4. { 'Height': 200, 'Weight': 90, 'Age': 35, 'Speed': 110, 'Tax Amount': 2500},
5. { 'Height': 210, 'Weight': 100, 'Age': 40, 'Speed': 120, 'Tax Amount': 3000}]

Distinguishing qualitative and quantitative data

As discussed above, qualitative data describes the quality or nature of something, such as color, shape, taste, or opinion etc. whereas quantitative data is the data that measures the quantity or amount of something, such as length, weight, speed, or frequency. Qualitative data can be further classified as nominal (categorical) or ordinal (ranked). Quantitative data can be further classified as discrete (countable) or continuous (measurable). The following methods are used to understand if data is qualitative or quantitative in nature.

dtype(): It is used to check the data types of the data frame.

Tutorial 1.11: To implement **dtype()** to check the datatypes of the different features or column in a data frame, as follows:

1. `import pandas as pd`
2. *# Create a data frame with qualitative and quantitative*

columns

```
3. df = pd.DataFrame({
4.     "age": [25, 30, 35], # a quantitative column
5.     "gender": ["female", "male", "male"], # a qualitative
   column
6.     "hair color": ["black", "brown", "white"], # a qualitati
   ve column
7.     "marital status": ["single", "married", "divorced"], #
   a qualitative column
8.     "salary": [5000, 6000, 7000], # a quantitative column
9.     "height": [6, 5.7, 5.5], # a quantitative column
10.    "weight": [60, 57, 55] # a quantitative column
11. })
12. # Print the data frame
13. print(df)
14. # Print the data types of each column using dtype()
15. print(df.dtypes)
```

Output:

```
1.  age  gender  hair color  marital status  salary  height  w
   eight
2.  0   25  female    black        single   5000    6.0    60
3.  1   30   male    brown        married   6000    5.7    57
4.  2   35   male    white        divorced   7000    5.5    55
5.  age                int64
6.  gender              object
7.  hair color          object
8.  marital status      object
9.  salary              int64
10. height              float64
11. weight              int64
```

describe(): You can also use the describe method from pandas to generate descriptive statistics for each column. This method will only show statistics for quantitative columns by default, such as mean, standard deviation, minimum, maximum, etc.

You need to specify **include='O'** as an argument to include qualitative columns. This will show statistics for qualitative columns, such as count, unique values, top values, and frequency. As you can see, the descriptive statistics for qualitative and quantitative columns are different, reflecting the nature of the data.

Tutorial 1.12: To implement **describe()** in the data frame used in *Tutorial 1.11* of **dtype()**, is as follows:

1. *# Print the descriptive statistics for quantitative columns*
2. `print(df.describe())`
3. *# Print the descriptive statistics for qualitative columns*
4. `print(df.describe(include='O'))`

Output:

1.	age	salary	height	weight
2. count	3.0	3.0	3.000000	3.000000
3. mean	30.0	6000.0	5.733333	57.333333
4. std	5.0	1000.0	0.251661	2.516611
5. min	25.0	5000.0	5.500000	55.000000
6. 25%	27.5	5500.0	5.600000	56.000000
7. 50%	30.0	6000.0	5.700000	57.000000
8. 75%	32.5	6500.0	5.850000	58.500000
9. max	35.0	7000.0	6.000000	60.000000
10.	gender	hair	color	marital status
11. count	3	3	3	
12. unique	2	3	3	
13. top	male	black	single	
14. freq	2	1	1	

value_counts(): To count unique values in a data frame, the **value_counts()** is used. It also displays the data type **dtype**. The **dtype** displays is the data type of the values in the series object returned by the **value_counts** method.

Tutorial 1.13: To implement **value_count()** to count unique value in a data frame as follows:

```
1. # To count the values in `gender` column
2. print(df['gender'].value_counts())
3. print("\n")
4. # To count the values in `age` column
5. print(df['age'].value_counts())
```

Output:

```
1. gender
2. male    2
3. female  1
4. Name: count, dtype: int64
5.
6. age
7. 25    1
8. 30    1
9. 35    1
```

In above *Tutorial 1.13* of **value_counts()**, the values are counts of each unique value in the gender column of the data frame, and the data type is **int64**, which means 64-bit integer.

is_numeric_dtype(), is_string_dtype(): These functions from the **pandas.api.types** module can help you determine if a column contains numeric or string (object) data.

Tutorial 1.14: To implement checking the numeric and string data type of a data frame column with **is_numeric_dtype()** and **is_string_dtype()** functions is as follows:

1. *# Import module for data type checking and inference.*
2. `import pandas.api.types as ptypes`
3. *# Checks if the column 'hair color' in df is of the string dtype and prints the result*
4. `print(f"Is string?: {ptypes.is_string_dtype(df['hair color'])}")`
5. *# Checks if the column 'weight' in df is of the numeric dtype and prints the result*
6. `print(f"Is numeric?: {ptypes.is_numeric_dtype(df['weight'])}")`
7. *# Checks if the column 'salary' in df is of the string dtype and prints the result*
8. `print(f"Is string?: {ptypes.is_string_dtype(df['salary'])}")`

Output:

1. Is string?: `True`
2. Is numeric?: `True`
3. Is string?: `False`

Also, in *Tutorial 1.14* we can use a for loop to check it for all the columns iteratively as follows:

1. *# Check the data types of each column using is_numeric_dtype() and is_string_dtype()*
2. `for col in df.columns:`
3. `print(f"{col}:")`
4. `print(f"Is numeric? {ptypes.is_numeric_dtype(df[col])}")`
5. `print(f"Is string? {ptypes.is_string_dtype(df[col])}")`
6. `print()`

info(): It describes the data frame with a column name, the number of not null values, and the data type of each column.

Tutorial 1.15: To implement **info()** to view the information

about a data frame is as follows:

```
1. df.info()
```

The output will display the summary consisting of column names, **non-null** count values, data types of each column, and many more.

```
1. RangeIndex: 3 entries, 0 to 2
```

```
2. Data columns (total 7 columns):
```

```
3. #   Column      Non-Null Count  Dtype
```

```
4. ---  ---
```

```
5. 0   age          3 non-null    int64
```

```
6. 1   gender        3 non-null    object
```

```
7. 2   hair color     3 non-null    object
```

```
8. 3   marital status 3 non-null    object
```

```
9. 4   salary         3 non-null    int64
```

```
10. 5   height         3 non-null    float64
```

```
11. 6   weight         3 non-null    int64
```

```
12. dtypes: float64(1), int64(3), object(3)
```

```
13. memory usage: 296.0+ bytes
```

head() and tail(): **head()** displays the data frame from the top, and the **tail()** displays it from the last or bottom.

Tutorial 1.16: To implement **head()** and **tail()** to view the top and bottom rows of data frame respectively.

head() displays the data frame with the first few rows. Inside the parenthesis of **head()**, we can define the number of rows we want to view. For example, to view the first ten rows of the data frame, write **head(10)**. Same with **tail()** also as follows:

```
1. # View the first few rows of the DataFrame
```

```
2. print(df.head())
```

```
3. # View first 2 rows of the DataFrame
```

```
4. df.head(1)
```

The output will display the data frame from the top, and **head(1)** will display the topmost row of data frame as follows:

```
1. age gender hair color marital status salary height weight
2. 0 25 female black single 5000 6.0 60
3. 1 30 male brown married 6000 5.7 57
4. 2 35 male white divorced 7000 5.5 55
5. age gender hair color marital status salary height weight
6. 0 25 female black single 5000 6.0 60
```

Tutorial 1.17: To implement **tail()** and display the bottom most rows of data frame is as follows:

```
1. # Import display function from IPython module to display the dataframe
2. from IPython.display import display
3. # View the last few rows of the DataFrame
4. display(df.tail())
5. # View last 2 rows of the DataFrame
6. display(df.tail(1))
```

In output **tail(1)** will display the bottommost row of data frame as follows:

```
1. age gender hair color marital status salary height weight
2. 0 25 female black single 5000 6.0 60
3. 1 30 male brown married 6000 5.7 57
4. 2 35 male white divorced 7000 5.5 55
5. age gender hair color marital status salary height weight
6. 2 35 male white divorced 7000 5.5 55
```

Other methods: Besides function described above there

are few other methods in Python that are useful to distinguish qualitative and quantitative data. They are as follows:

- **select_dtypes(include='__')**: It is used to select columns with specific datatype that is, number and object.

1. *# Select and display DataFrame with only numeric values*

2. `display(df.select_dtypes(include='number'))`

The output will display only the numeric column of the data frame as follows:

1. `age salary height weight`

2. `0 25 5000 6.0 60`

3. `1 30 6000 5.7 57`

4. `2 35 7000 5.5 55`

To select only object data types, an object is used. In pandas, objects are the column containing strings or mixed types of data. It is the default data type for columns that have text or arbitrary Python objects as follows:

1. *# Select and display DataFrame with only object values in the same cell(display() displays DataFrame in same cell)*

2. `display(df.select_dtypes(include='object'))`

Output will include only object type columns as follows:

1. `gender hair color marital status`

2. `0 female black single`

3. `1 male brown married`

4. `2 male white divorced`

- **groupby():** `groupby()` is used to group based on column as follows:

1. *# Group by gender*

2. `df.groupby("gender")`

After grouping based on the column name, it can be used to display the summary statistics of a grouped data frame. **describe()** method on the **groupby()** object can be used to find summary statistics of each group as follows:

```
1. # Describe dataframe summary statistics by gender
2. df.groupby('gender').describe()
```

Further to print the count of each group, you can use the **size()** or **count()** method on the **groupby()** object as follows:

```
1. # Print count of group object with size
2. print(df.groupby('gender').size())
3. # Print count of group object with count
4. print(df.groupby('gender').count())
```

Output:

```
1. gender
2. female    1
3. male      2
4. dtype: int64
5. age  hair color  marital status  salary  height  weight
6. gender
7. female    1    1              1      1      1      1
8. male      2    2              2      2      2      2
```

- **groupby().sum():** **groupby().sum()** groups data and then display sum in each group as follows:

```
1. # Group by gender and hair color and calculate the sum of each group
2. df.groupby(["gender", "hair color"]).sum()
```

- **columns:** **columns** display column names. Sometimes, through descriptive columns names, types of data can

be distinguished. So, displaying column name can be useful as follows:

1. *# Displays all column names.*

2. `df.columns`

- **type():** **type()** is used to display the type of a variable. It can be used to determine the type of a single variable as follows:

1. *# Declare variable*

2. `x = 42`

3. `y = "Hello"`

4. `z = [1, 2, 3]`

5. *# Print data types*

6. `print(type(x))`

7. `print(type(y))`

8. `print(type(z))`

Tutorial 1.18: To implement **read_json()**, to read and view nobel prize dataset in JSON format.

Let us load a *nobel prizdataset*² and see what kind of data it contains. The *Tutorial 1.18* flattens nested **JSON** data structures into a data frame as follows:

1. `import pandas as pd`

2. *# Read the json file from the direcotory*

3. `json_df = pd.read_json("/workspaces/ImplementingStatisticsWithPython/data/chapter1/prize.json")`

4. *# Convert the json data into a dataframe*

5. `data = json_df["prizes"]`

6. `prize_df = pd.json_normalize(data)`

7. *# Display the dataframe*

8. `prize_df`

To see what type of data **prize_df** contains use **info()** and **head()**, is as follows:

1. `prize_df.info()`

```
2. prize_df.head()
```

Alternatively, to *Tutorial 1.18*, the *nobel prize dataset*³ can be accessed directly by sending the request as shown in the following code:

```
1. import pandas as pd
2. # Send HTTP requests using Python
3. import requests
4. # Get the json data from the url
5. response = requests.get("https://api.nobelprize.org/v1/
    prize.json")
6. data = response.json()
7. # Convert the json data into a dataframe
8. prize_json_df = pd.json_normalize(data, record_path="p
    rizes")
9. prize_json_df
```

Tutorial 1.19: To implement **read_csv()**, to read and view nobel prize dataset in CSV format is as follows:

```
1. import pandas as pd
2. # Read the json file from the direcotory
3. prize_csv_df = pd.read_csv("/workspaces/Implementing
    StatisticsWithPython/data/chapter1/prize.csv")
4. # Display the dataframe
5. prize_csv_df
```

Tutorial 1.20: To implement use of NumPy and to read diabetes dataset in CSV files.

Most common ways are using **numpy.loadtxt()** and using **numpy.genfromtxt()**. **numpy.loadtxt()** assumes that the file has no missing values, no comments, and no headers and uses whitespace as the delimiter by default. We can change the delimiter to a comma by passing **delimiter = ','** as a parameter. Here, the CSV file has one header row, which is a string, so we use **skiprows = 1** this skips the

first row of the CSV file and loads the rest of the data as a NumPy array as follows:

```
1. import numpy as np
2. arr = np.loadtxt('/workspaces/ImplementingStatisticsWithPython/data/chapter1/diabetes.csv', delimiter=',', skiprows=1)
3. print(arr)
```

The **numpy.genfromtxt()** function can handle missing values, comments, headers, and various delimiters. We can use the **missing_values** parameter to specify which values to treat as missing. We can use the **comments** parameter to specify which character indicates a comment line, such as # or %. For example, if you have a CSV file named **diabetes.csv** that looks as follows:

```
1. import numpy as np
2. arr = np.genfromtxt('/workspaces/ImplementingStatisticsWithPython/data/chapter1/diabetes.csv', delimiter=',', names=True, missing_values='?', dtype=None)
3. print(arr)
```

Univariate, bivariate, and multivariate data

Univariate, bivariate, and multivariate data are terms used in statistics to describe the number of variables and their relationships within a dataset. Where univariate means one, bivariate means two and multivariate is more than two. These concepts are fundamental to statistical analysis and play a crucial role in various fields, from social sciences to natural sciences, engineering and beyond.

Univariate data and univariate analysis

Univariate analysis involves observing only one variable or attribute. For example, height of students in a class, color of cars in a parking lot, or salary of employees in a

company are all univariate data. Univariate analysis analyzes only one variable column or attribute at a time. For example, analyzing only the patient height column at a time or the person's salary column.

Tutorial 1.21: To implement univariate data and univariate analysis by selecting a column or variable or attribute from the CSV dataset and compute its mean, standard deviation, frequency or distribution with other information using **describe()** as follows:

```
1. import pandas as pd
2. from IPython.display import display
3. diabeties_df = pd.read_csv('/workspaces/ImplementingS
   tatisticsWithPython/data/chapter1/diabetes.csv')
4. #To view all the column names
5. print(diabeties_df.columns)
6. # Select the column Glucose column as a DataFrame fro
   m diabeties_df DataFrame
7. display(diabeties_df[['Glucose']])
8. # describe() gives the mean,standard deviation
9. print(diabeties_df[['Glucose']].describe())
```

In the above, we selected only the **Glucose** column of the **diabeties_df** and analyzed only that column. This kind of single-column analysis is univariate analysis.

Tutorial 1.22: To further implement computation of median, mode range, frequency or distribution of variables in continuation with *Tutorial 1.21*, is as follows:

```
1. # Use mode() for computing most frequest value i.e, mo
   de
2. print(diabeties_df[['Glucose']].mode())
3. # To get range simply subtract DataFrame maximum val
   ue by the DataFrame minimum value. Use df.max() and
   df.min() for maximum and minimum value
4. mode_range = diabeties_df[['Glucose']].max() - diabeties
```

```

    _df[['Glucose']].min()
5. print(mode_range)
6. # For frequency or distribution of variables use value_counts()
7. diabeties_df[['Glucose']].value_counts()

```

Bivariate data

Bivariate data consists of observing two variables or attributes for each individual or unit. For example, if you wanted to study the relationship between the age and height of students in a class, you would collect the age and height of each student. Age and height are two variables or attributes, and each student is an individual or unit. Bivariate analysis analyzes how two different variables, columns, or attributes are related. For example, the correlation between people's height and weight or between hours worked and monthly salary.

Tutorial 1.23: To implement bivariate data and bivariate analysis by selecting two columns or variables or attributes from the CSV dataset and to describe them, as follows:

```

1. import pandas as pd
2. from IPython.display import display
3. diabeties_df = pd.read_csv('/workspaces/ImplementingStatisticsWithPython/data/chapter1/diabetes.csv')
4. #To view all the column names
5. print(diabeties_df.columns)
6. # Select two column Glucose column as a DataFrame from diabeties_df DataFrame
7. display(diabeties_df[['Glucose','Age']])
8. # describe() gives the mean,standard deviation
9. print(diabeties_df[['Glucose','Age']].describe())
10. # Use mode() for computing most frequest value i.e, mo

```

de

```
11. print(diabities_df[['Glucose']].mode())
12. # To get range simply subtract DataFrame maximum value by the DataFrame minimum value. Use df.max() and df.min() for maximum and minimum value
13. mode_range = diabities_df[['Glucose']].max() - diabities_df[['Glucose']].min()
14. print(mode_range)
15. # For frequency or distribution of variables use value_counts()
16. diabities_df[['Glucose']].value_counts()
```

Here, we compared two columns, **glucose** and **age** in **diabities_df** data frame, which involved multiple data frame columns making it bivariate analysis.

Alternatively two or more columns can be accessed using **loc[row_start:row:stop,column_start:column:stop]** or also through column index via slicing by using **iloc[row_start:row:stop,column_start:column:stop]** as follows:

```
1. # Using loc
2. diabities_df.loc[:, ['Glucose','Age']]
3. # Using iloc, column index and slicing
4. diabities_df.iloc[:,0:2]
```

Further, to compute the correlation between two variables or two columns, such as **glucose** and **age**, we can use columns along with **corr()** as follows:

```
1. diabities_df['Glucose'].corr(diabities_df['Age'])
```

Correlation is a statistical measure that indicates how two variables are related to each other. A positive correlation means that the variables increase or decrease together, while a negative correlation means that the variables move in opposite directions. A correlation value close to zero

means that there is no linear relationship between the variables.

In the context of `diabetes_df['Glucose'].corr(diabetes_df['Age'])`, the random positive correlation value of **0.26** means that there is a weak positive correlation between glucose level and age in the diabetes dataset. This implies that older people tend to have higher glucose levels than younger people but the relationship is not very strong or consistent. Correlation can be computed using different methods such as **pearson**, **kendall**, or **spearman** then specify `method = '__'` in `corr()` as follows:

```
1. diabetes_df['Glucose'].corr(diabetes_df['Age'], method='kendall')
```

Multivariate data

Multivariate data consists of observing three or more variables or attributes for each individual or unit. For example, if you want to study the relationship between the age, gender, and income of customers in a store, you would collect this data for each customer. Age, gender, and income are the three variables or attributes, and each customer is an individual or unit. In this case, the data you collect will be multivariate data because it requires observations on three variables or attributes for each individual or unit. For example, the correlation between age, gender, and sales in a store or between temperature, humidity, and air quality in a city.

Tutorial 1.24: To implement multivariate data and multivariate analysis by selecting multiple columns or variables or attributes from the CSV dataset and describe them, as follows:

```
1. import pandas as pd
2. from IPython.display import display
```

```

3. diabetics_df = pd.read_csv('/workspaces/ImplementingS
   tatisticsWithPython/data/chapter1/diabetes.csv')
4. #To view all the column names
5. print(diabetics_df.columns)
6. # Select the column Glucose column as a DataFrame fro
   m diabetics_df DataFrame
7. display(diabetics_df[['Glucose','BMI', 'Age', 'Outcome']])
8. # describe() gives the mean,standard deviation
9. print(diabetics_df[['Glucose','BMI', 'Age', 'Outcome']].de
   scribe())

```

Alternatively, multivariate analysis can be performed by describing the whole data frame as follows:

```

1. # describe() gives the mean,standard deviation
2. print(diabetics_df.describe())
3. # Use mode() for computing most frequent value i.e, mo
   de
4. print(diabetics_df.mode())
5. # To get range simply subtract DataFrame maximum val
   ue by the DataFrame minimum value. Use df.max() and
   df.min() for maximum and minimum value
6. mode_range = diabetics_df.max() - diabetics_df.min()
7. print(mode_range)
8. # For frequency or distribution of variables use value_co
   unts()
9. diabetics_df.value_counts()

```

Further, to compute the correlation between all the variables in the data frame, use **corr()** after the data frame variable name as follows:

```

1. diabetics_df.corr()

```

You can also apply various multivariate analysis techniques, as follows:

- **Principal Component Analysis (PCA):** It transforms

high-dimensional data into a smaller set of uncorrelated variables (principal components) that capture the most variance, thereby simplifying the dataset while retaining essential information. It makes easier to visualize, interpret, and model multivariate relationships

- **Library:** Scikit-learn
- **Method:** **PCA(n_components=__)**
- **Multivariate regression:** This is used to analyze the relationship between multiple dependent and independent variables.
 - **Library:** Statsmodels
 - **Method:** **statsmodels.api.OLS** for ordinary least squares regression. It allows you to perform multivariate linear regression and analyze the relationship between multiple dependent and independent variables. Regression can also be performed using scikit-learn's **LinearRegression()**, **LogisticRegression()**, and many more.
- **Cluster analysis:** This is used to group similar data points together based on their characteristics.
 - **Library:** Scikit-learn
 - **Method:** **sklearn.cluster.KMeans** for K-means clustering. It allows you to group similar data points together based on their characteristics. And many more.
- **Factor analysis:** This is used to identify underlying latent variables that explain the observed variance.
 - **Library:** FactorAnalyzer
 - **Method:** **FactorAnalyzer** for factor analysis. It allows you to perform **Exploratory Factor**

Analysis (EFA) to identify underlying latent variables that explain the observed variance.

- **Canonical Correlation Analysis (CCA):** To explore the relationship between two sets of variables.
 - **Library:** Scikit-learn
 - **Method:** **sklearn.cross_decomposition** and CCA allows you to explore the relationship between two sets of variables and find linear combinations that maximize the correlation between the two sets.

Tutorial 1.25: To implement **Principal Component Analysis (PCA)** for dimensionality reduction is as follows:

1. `import pandas as pd`
2. `# Import principal component analysis`
3. `from sklearn.decomposition import PCA`
4. `# Scales data between 0 and 1`
5. `from sklearn.preprocessing import StandardScaler`
6. `# Import matplotlib to plot visualization`
7. `import matplotlib.pyplot as plt`
8. `# Step 1: Load your dataset into a DataFrame`
9. `# Assuming you have your dataset stored in a CSV file called "data.csv", load it into a Pandas DataFrame.`
10. `data = pd.read_csv("/workspaces/ImplementingStatisticsWithPython/data/chapter1/diabetes.csv")`
11. `# Step 2: Separate the features and the outcome variable (if applicable)`
12. `# If the "Outcome" column represents the dependent variable and not a feature, you should separate it from the features.`
13. `# If it's not the case, you can skip this step.`
14. `X = data.drop("Outcome", axis=1) # Features`


```
15. y = data["Outcome"] # Outcome (if applicable)
16. # Step 3: Standardize the features
17. # PCA is sensitive to the scale of features, so it's crucial
    to standardize them to have zero mean and unit variance.
18. scaler = StandardScaler()
19. X_scaled = scaler.fit_transform(X)
20. # Step 4: Apply PCA for dimensionality reduction
21. # Create a PCA instance and specify the number of components
    you want to retain.
22. # If you want to reduce the dataset to a certain number
    of dimensions (e.g., 2 or 3), set the 'n_components' accordingly.
23. pca = PCA(n_components=2) # Reduce to 2 principal components
24. X_pca = pca.fit_transform(X_scaled)
25. # Step 5: Explained Variance Ratio
26. # The explained variance ratio gives us an idea of how much
    information each principal component captures.
27. explained_variance_ratio = pca.explained_variance_ratio_
28. # Step 6: Visualize the Explained Variance Ratio
29. plt.bar(range(len(explained_variance_ratio)), explained_
    variance_ratio)
30. plt.xlabel("Principal Component")
31. plt.ylabel("Explained Variance Ratio")
32. plt.title("Explained Variance Ratio for Each Principal Component")
33. # Show the figure
34. plt.savefig('skew_negative.jpg',dpi=600,bbox_inches='tight')
```

```
35. plt.show()
```

PCA reduces the dimensions but it also results in some loss of information as we only retain the most important components. Here, the original 8-dimensional diabetes data set has been transformed into a new 2-dimensional data set. The two new columns represent the first and second principal components, which are linear combinations of the original features. These principal components capture the most significant variation in the data.

The columns of the data set pregnancies, glucose, blood pressure, skin thickness, insulin, BMI, diabetes pedigree function, and age are reduced to 2 principal components because we specify **n_components=2** as shown in [Figure 1.1](#).

Output:

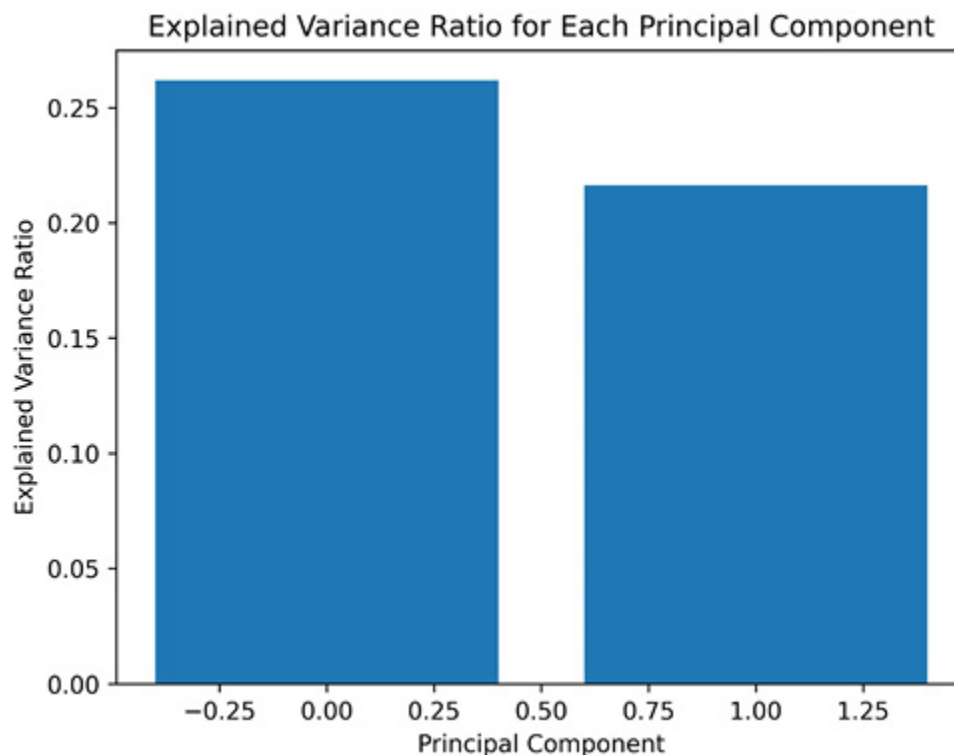


Figure 1.1: Explained variance ratio for each principal component

Following is what you can infer from these explained

variance ratios in this **diabetes dataset**:

- **The First Principal Component (PC1):** With an explained variance of 0.27, PC1 captures the largest portion of the data's variability. It represents the direction in the data space along which the data points exhibit the most significant variation. PC1 is the principal component that explains the most significant patterns in the data.
- **The Second Principal Component (PC2):** With an explained variance of 0.23, PC2 captures the second-largest portion of the data's variability. PC2 is orthogonal (uncorrelated) to PC1, meaning it represents a different direction in the data space from PC1. PC2 captures additional patterns that are not explained by PC1 and provides complementary information. PC1 and PC2 account for approximately 50% ($0.27 + 0.23$) of the total variance.

You can do similar with NumPy and JSON. Also, you can create different types of plots and charts for data analysis using Matplotlib and Seaborn libraries

Data sources, methods, populations, and samples

Data sources provide information for analysis from surveys, databases, or experiments. Collection methods determine how data is gathered, through interviews, questionnaires, or observations. Population is the entire group being studied, while samples are representative subsets used to draw conclusions with less analysis.

Data source

Data can be primary and secondary. It can be of two types, that is, statistical sources like surveys, census, experiments, and statistical reports and non-statistical

sources like business transactions, social media posts, weblogs, data from wearables and sensors, or personal records.

Tutorial 1.26: To implement reading data from different sources and view statistical and non-statistical data is as follows:

1. `import pandas as pd`
2. *# To import urllib library for opening and reading URLs*
3. `import urllib.request`
4. *# To access CSV file replace file name*
5. `df = pd.read_csv('url_to_csv_file.csv')`

To access or read data from different sources, pandas provides **read_csv()** and **read_json()** and **loadtxt()**, **genfromtxt()** in NumPy and many others. The URL can also be used like <https://api.nobelprize.org/v1/prize.json>, but it should be accessible. Most data server would need authentication to access the server.

To read JSON files replace file name in the script as follows:

1. *# To access JSON data replace file name*
2. `df = pd.read_json('your_file_name.json')`

To read XML file from a server with NumPy, you can use the **np.loadtxt()** function and pass as an argument a file object created using the **urllib.request.urlopen()** function from the **urllib.request** module. You must also specify the delimiter parameter as `<` or `>` to separate XML tags from the data values. To read XML file, replace files names with appropriate one in the script as follows:

1. *# To access and read the XML file using URL*
2. `file = urllib.request.urlopen('your_url_to_accessible_xml_file.xml')`
3. *# To open the XML file from the URL and store it in a file*

object

```
4. arr = np.loadtxt(file, delimiter='<')  
5. print(arr)
```

Collection methods

Collection methods are surveys, interviews, observations, focus groups, experiments, and secondary data analysis. It can be quantitative, based on numerical data and statistical analysis, or qualitative, based on words, images, actions, and interpretive analysis. Also, sometimes mixed methods, which combine qualitative and quantitative, can be used.

Population and sample

Population is entire group of people, items, or elements you want to study or draw conclusions about. For example, if you want to know the average score of all students in a school, the population is all students. The sample is a subset of the population from which you select and collect data. For example, 20 randomly chosen students from this school are a population sample.

Let us see an example of selecting population and sample using random modules. The **random** module **rand()** function can be utilized to randomly choose unstructured and semi-structured datasets or files. This approach ensures that each data point has an equal probability of being included in the sample, thereby minimizing selection bias and ensuring the sample's representativeness of the broader population.

Tutorial 1.27: To implement **rand()** to select items from the population, is as follows:

```
1. import random  
2. # Define population and sample size  
3. population = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]  
4. sample_size = 3
```

5. *# Randomly select a sample from the population*
6. `sample = random.sample(population, sample_size)`
7. `print("Sample:", sample)`

Tutorial 1.28: To implement **rand()** to select items from the patient registry data, is as follows:

1. `import pandas as pd`
2. `from IPython.display import display`
3. *# Import module to generate random numbers*
4. `import random`
5. *# Read CSV file and save as dataframe*
6. `diabities_df = pd.read_csv('/workspaces/ImplementingStatisticsWithPython/data/chapter1/diabetes.csv')`
7. *# Define the sample size*
8. `sample_size = 5`
9. *# Get the number of rows in the DataFrame*
10. `num_rows = diabities_df.shape[0]`
11. *# Generate random indices for selecting rows*
12. `random_indices = random.sample(range(num_rows), sample_size)`
13. *# Select the rows using the random indices*
14. `sample_diabities_df = diabities_df.iloc[random_indices]`
15. `display(sample_diabities_df)`

While random sampling methods like **rand()** help select a representative subset from a broader population, functions such as **train_test_split()** play a pivotal role in organizing this subset into training and testing sets, particularly in supervised learning. By systematically dividing data into dependent and independent variables and ensuring that these splits are both representative and reproducible, **train_test_split()** facilitates the development of models that perform reliably on unseen data.

Tutorial 1.29: To implement **train_test_split()** to select

items from the patient registry data population, is as follows:

```
1. # Import sklearn train_test_split
2. from sklearn.model_selection import train_test_split
3. # Define population and test size
4. population = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
5. test_size = 0.2 # Proportion of data to reserve for testing
6. # Split the population into training and testing sets
7. train_set, test_set = train_test_split(population, test_size=test_size, random_state=42)
8. # Display the split
9. print("Training Set:", train_set)
10. print("Testing Set:", test_set)
```

Output:

```
1. Training Set: [6, 1, 8, 3, 10, 5, 4, 7]
2. Testing Set: [9, 2]
```

Data preparation tasks

Data preparation tasks are the early steps carried out upon having access to the data. It involves checking the quality of data, cleaning of data, data wrangling, and its manipulation described in detail.

Data quality

Data quality indicates how suitable, accurate, useful, complete, reliable, and consistent the data is for its intended use. Verifying data quality is an important step in analysis and preprocessing.

Tutorial 1.30: To implement checking the data quality of CSV file data frame, is as follows:

- Check missing values with **isna()** or **isnull()**
- Check summary with **describe()** or **info()**
- Check shape with **shape**, size with **size**, and memory usage with **memory_usage()**
- Check duplicates with **duplicated()** and remove duplicate with **drop_duplicates()**

Based on this instruction, let us see the implementation as follows:

```
1. import pandas as pd
2. diabeties_df = pd.read_csv('/workspaces/ImplementingS
   tatisticsWithPython/data/chapter1/diabetes.csv')
3. # Check for missing values using isna() or isnull()
4. print(diabeties_df.isna().sum())
5. #Describe the dataframe with describe() or info()
6. print(diabeties_df.describe())
7. # Check for the shape,size and memory usage
8. print(f'Shape: {diabeties_df.shape} Size: {diabeties_df.siz
   e} Memory Usage: {diabeties_df.memory_usage()}')
9. # Check for the duplicates using duplicated() and drop t
   hem if necessary using drop_duplicates()
10. print(diabeties_df.duplicated())
```

Now, we use synthetic transaction narrative data containing unstructured information about the nature of the transaction.

Tutorial 1.30: To implement viewing the text information in the text files (synthetic transaction narrative files), is as follows:

```
1. import pandas as pd
2. import numpy as np
3. # To import glob library for finding files and directories u
   sing patterns
4. import glob
```



```

5. # To assign the path of the directory containing the text
   files to a variable
6. path = "/workspaces/ImplementingStatisticsWithPython
   /data/chapter1/TransactionNarrative"
7. # To find all the files in the directory that have a .txt ext
   ension and store them in a list
8. files = glob.glob(path + "/*.txt")
9. # To loop through each file in the list
10. for file in files:
11.     # To open each file in read mode with utf-
       8 encoding and assign it to a file object
12.     with open(file, "r", encoding="utf-8") as f:
13.         print(f.read())

```

Output:

1. Date: 2023-08-01
2. Merchant: VideoStream Plus
3. Amount: \$9.99
4. Description: Monthly renewal of VideoStream Plus subs
cription.
5. Your subscription to VideoStream Plus has been success
fully renewed for \$9.99.

Tutorial 1.31: To implement checking the data quality of multiple **.txt** files (synthetic transaction narrative files) that contains text information as shown in *Tutorial 1.30* output and to check the quality of information in them, we use **file_size**, **line_count**, **missing_field**, as follows:

```

1. import os
2. import glob
3. def check_file_quality(content):
4.     # Check for presence of required fields
5.     required_fields = ['Date:', 'Merchant:', 'Amount:', 'De
       scription:']

```

```
6. missing_fields = [field for field in required_fields if field
    not in content]
7. # Calculate file size
8. file_size = len(content.encode('utf-8'))
9. # Count lines in the content
10. line_count = content.count('\n') + 1
11. # Return quality assessment
12. quality_assessment = {
13.     "file_name": file,
14.     "file_size_bytes": file_size,
15.     "line_count": line_count,
16.     "missing_fields": missing_fields
17. }
18. return quality_assessment
19. # To assign the path of the directory containing the text
    files to a variable
20. path = "/workspaces/ImplementingStatisticsWithPython
    /data/chapter1/TransactionNarrative"
21. # To find all the files in the directory that have a .txt ext
    ension and store them in a list
22. files = glob.glob(path + "/*.txt")
23. # To loop through each file in the list
24. for file in files:
25.     with open(file, "r", encoding="utf-8") as f:
26.         content = f.read()
27.         print(content)
28.         quality_result = check_file_quality(content)
29.         print(f"\nQuality Assessment for {quality_result['fil
    e_name']}:")
30.         print(f"File Size: {quality_result['file_size_bytes']} b
    ytes")
```

```

31.     print(f"Line Count: {quality_result['line_count']} lin
      es")
32.     if quality_result['missing_fields']:
33.         print("Missing Fields:", ', '.join(quality_result['mi
      ssing_fields']))
34.     else:
35.         print("All required fields present.")
36.     print("=" * 40)

```

Output (Only one transaction narrative output is shown):

```

1. Date: 2023-08-01
2. Merchant: VideoStream Plus
3. Amount: $9.99
4. Description: Monthly renewal of VideoStream Plus subs
  cription.
5.
6. Your subscription to VideoStream Plus has been success
  fully renewed for $9.99.
7.
8.
9. Quality Assessment for /workspaces/ImplementingStati
  sticsWithPython/data/chapter1/TransactionNarrative/3.
  txt:
10. File Size: 201 bytes
11. Line Count: 7 lines
12. All required fields present.
13. =====
    =====

```

Cleaning

Data cleansing involves identifying and resolving inconsistencies and errors in raw data sets to improve data quality. High-quality data is critical to gaining accurate and

meaningful insights. Data cleansing also include data handling. Different ways for data cleaning or data handling are described below.

Missing values

Missing values refer to data points or observations with incomplete or absent information. For example, in a survey, if people do not answer a certain question, the related entries will be empty. Appropriate methods, like imputation or exclusion, are used to address them. If there are missing values then one way is to drop missing value as shown in *Tutorial 1.32*.

Tutorial 1.32: To implement finding the missing value and dropping them.

Let us check **prize_csv_df** data frame for null values and drop the null ones, as follows:

```
1. import pandas as pd
2. from IPython.display import display
3. # Read the prize csv file from the direcotory
4. prize_csv_df = pd.read_csv("/workspaces/Implementing
    StatisticsWithPython/data/chapter1/prize.csv")
5. # Display the dataframe null values count
6. print(prize_csv_df.isna().sum())
```

Output:

```
1. year                374
2. category            374
3. overallMotivation   980
4. laureates__id       49
5. laureates__firstname 50
6. laureates__surname  82
7. laureates__motivation 49
8. laureates__share     49
```

Since **prize_csv_df** have null values, let us drop them and view the count of null values after drop as follows:

1. `print("\n \n **** After dropping the null values in prize_csv_df****")`
2. `after_dropping_null_prize_df = prize_csv_df.dropna()`
3. `print(after_dropping_null_prize_df.isna().sum())`

Finally, after applying the above code, the output will be as follows:

1. `**** After dropping the null values in prize_csv_df****`
2. year 0
3. category 0
4. overallMotivation 0
5. laureates__id 0
6. laureates__firstname 0
7. laureates__surname 0
8. laureates__motivation 0
9. laureates__share 0
10. dtype: int64

This shows there are now zero null values in all the column.

Imputation

Imputation means to place a substitute value in place of the missing values. Like constant value imputation, mean imputation, mode imputation.

Tutorial 1.33: To implement imputing the mean value of the column **laureates__share**.

Mean imputation only imputes the mean value of numeric data types as **fillna()** expects scalar, so we cannot use the **mean()** method to fill missing values in object columns.

1. `import pandas as pd`
2. `from IPython.display import display`
3. `# Read the prize csv file from the direcotory`

```

4. prize_csv_df = pd.read_csv("/workspaces/Implementing
   StatisticsWithPython/data/chapter1/prize.csv")
5. # View the number of null values in original DataFrame
6. print("Null Value Before",prize_csv_df['laureates__share']
   .isna().sum())
7. # Calculate the mean of each column
8. prize_col_mean = prize_csv_df['laureates__share'].mean
   ()
9. # Fill missing values with column mean, inplace = True
   will replace the original DataFrame
10. prize_csv_df['laureates__share'].fillna(value=prize_col_
   mean, inplace=True)
11. # View the number of null values in the new DataFrame
12. print("Null Value After",prize_csv_df['laureates__share']
   .isna().sum())

```

Output:

1. Null Value Before 49
2. Null Value After 0

Also to fill missing values in object columns, you have to use a different strategy, such as using a constant value i.e, **df[column_name].fillna(' ')**, a mode value, or a custom function..

Tutorial 1.34: To implement imputing the mode value in the object data type column.

```

1. import pandas as pd
2. from IPython.display import display
3. # Read the prize csv file from the direcotory
4. prize_csv_df = pd.read_csv("/workspaces/Implementing
   StatisticsWithPython/data/chapter1/prize.csv")
5. # Display the original DataFrame null values in object d
   ata type columns
6. print(prize_csv_df.isna().sum())

```

```

7. # Select the object columns
8. object_cols = prize_csv_df.select_dtypes(include='object').columns
9. # Calculate the mode of each object data type column
10. col_mode = prize_csv_df[object_cols].mode().iloc[0]
11. # Fill missing values with the mode of each object data type column
12. prize_csv_df[object_cols] = prize_csv_df[object_cols].fillna(col_mode)
13. # Display the DataFrame column after filling null values in object data type columns
14. print(prize_csv_df.isna().sum())

```

Output:

```

1. year          374
2. category      374
3. overallMotivation  980
4. laureates__id    49
5. laureates__firstname  50
6. laureates__surname   82
7. laureates__motivation  49
8. laureates__share    49
9. dtype: int64
10. year          374
11. category      0
12. overallMotivation  0
13. laureates__id    49
14. laureates__firstname  0
15. laureates__surname   0
16. laureates__motivation  0
17. laureates__share    49
18. dtype: int64

```

Duplicates

Data may be duplicated or contains duplicate value. The duplicacy will affect the final statistical result. Hence, to prevent duplicacy, identifying and removing duplicates is necessary step as explained in this section. Best way to handle duplicate is to identify and remove duplicates.

Tutorial 1.35: To implement identifying and removing duplicate rows in data frame with **`duplicated()`**, as follows:

1. *# Identify duplicate rows and display their index*
2. `print(prize_csv_df.duplicated().index[prize_csv_df.duplicated()])`

Since, there is no duplicate the output is empty it displays indexes of duplicates as follows:

1. `Index([], dtype='int64')`

Also, you can find the duplicate values in a specific column by using the following code:

1. `prize_csv_df.duplicated(subset=['name_of_the_column'])`

To remove duplicates, **`drop()`** method can be used, syntax will be **`dataframe.drop(labels, axis='columns', inplace=False)`**. Drop can be applied to row and index using label and index values as follows:

1. `import pandas as pd`
2. *# Create a sample dataframe*
3. `people_df = pd.DataFrame({'name': ['Alice', 'Bob', 'Charlie'], 'age': [25, 30, 35], 'gender': ['F', 'M', 'M']})`
4. *# Print the original dataframe*
5. `print("original dataframe \n", people_df)`
6. *# Drop the 'gender' column and return a new dataframe*
7. `new_df = people_df.drop('gender', axis='columns')`
8. *# Print the new dataframe*
9. `print("dataframe after drop \n", new_df)`

Output:

```
1. original dataframe
2.      name age gender
3. 0  Alice  25    F
4. 1   Bob  30    M
5. 2 Charlie 35    M
6. dataframe after drop
7.      name age
8. 0  Alice  25
9. 1   Bob  30
10. 2 Charlie 35
```

Outliers

Outliers are data points that are very different from the other data points. They can be much higher or lower than the standard range of values. For example, if the heights of ten people in centimeters are measured, the values might be as follows:

160, 165, 170, 175, 180, 185, 190, 195, 200, 1500.

Most of the heights are alike but the last measurement is much larger than the others. This data point is an outlier because it is not like the rest of the data. The best way to handle outliers is to identify outliers and then correct, resolve, or leave as needed. Ways to identify outliers are to compute mean, standard deviation, and quantile (a common approach is to compute interquartile range). Another way to identify outliers is by computing the z-score of the data points and then considering points beyond the threshold values as outliers.

Tutorial 1.36: To implement identifying outliers in a data frame with `zscore`.

Z-score measures how many standard deviations a value is from the mean. In the following code, **z_score** identifies

outliers in the laureates' share column:

```
1. import pandas as pd
2. import numpy as np
3. # Read the prize csv file from the direcotory
4. prize_csv_df = pd.read_csv("/workspaces/Implementing
   StatisticsWithPython/data/chapter1/prize.csv")
5. # Calculate mean, standard deviation and Z-
   scores for the column
6. z_scores = np.abs((prize_csv_df['laureates__share'] - pri
   ze_csv_df['laureates__share'].mean()) / prize_csv_df['lau
   reates__share'].std())
7. # Define a threshold for outliers (e.g., 4)
8. threshold = 2
9. # Display the row index of the outliers
10. print(prize_csv_df.index[z_scores > threshold])
```

Output:

```
1. Index([ 17,  18,  22,  23,  34,  35,  48,  49,  54,  55
   ,
   62,  63,
2.      73,  74,  86,  87,  97,  98, 111, 112, 144, 145,
   146, 147,
3.      168, 169, 180, 181, 183, 184, 215, 216, 242,
   243,
   249, 250,
4.      255, 256, 277, 278, 302, 303, 393, 394, 425,
   426,
   467, 468,
5.      471, 472, 474, 475, 501, 502, 514, 515, 556,
   557,
   563, 564,
6.      607, 608, 635, 636, 645, 646, 683, 684, 760,
   761,
```

```
764, 765,  
7. 1022, 1023],  
8. dtype='int64')
```

The output shows the row index of the outliers in the laureates' share column of the **prize.csv** file. Outliers are values that are unusually high or low compared to the rest of the data. The code uses a z-score to measure how many standard deviations a value is from the mean of the column. A higher z-score means a more extreme value. The code defines a threshold of two, which means that any value with a z-score greater than two is considered an outlier.

Additionally, preparing data, cleaning it, manipulating it, and doing data wrangling includes the following:

- Checking typos and spelling errors. Python provides libraries like **PySpellChecker**, **NLTK**, **TextBlob**, or **Enchant** to check typos and spelling errors.
- Data transformation is a change from one form to another desired form. It involves aggregation, conversion, normalization, and many more, they are covered in detail in *Chapter 2, Exploratory Data Analysis*.
- Handling inconsistencies which involve identifying conflicting information and resolving them. For example, the body temperature is listed as 1400 Celsius which is not correct.
- Standardize format and units of measurements to ensure consistency.
- Further data integrity and validation ensures that data is unchanged, not altered or corrupted. Data validation verifies that the data to be used is correct (use techniques like validation rules, manual review).

Wrangling and manipulation

It means making raw data usable through cleaning, transformation, or other ways. It involves cleaning, organizing, merging, filtering, sorting, aggregating, and reshaping data. Helping you analyze, organize, and improve your data for informed insights and decisions.

The various useful data wrangling and manipulation method in Python are as following:

- **Cleaning:** Some of the methods used to clean the data, along with their syntax, are as follows:

`df.dropna()`: Removing missing values.

`df.fillna()`: Filling missing values.

`df.replace()`: Replacing values.

`df.drop_duplicates()`: Removing duplicate.

`df.drop()`: Removing specific rows **or** columns.

`df.rename()`: Renaming columns.

`df.astype()`: Changing data types.

- **Transformation:** Some of the methods used for data transformation, together with their syntax, are as follows:

`df.apply()`: Applying a function.

`df.groupby()`: Grouping data.

`df.pivot_table()`: Creating pivot tables to summarize.

`df.melt()`: Unpivoting **or** melting data.

`df.sort_values()`: Sorting rows.

`df.join()`, `df.merge()`: Combining data.

- **Aggregation:** Some methods used for data aggregation, together with their syntax, are as follows:

`df.groupby().agg()`: Aggregate data using specified functions.

`df.groupby().size()`, `df.groupby().count()`, `df.groupby().mean()`: Calculating common aggregation metrics.

- **Reshape:** Some methods used for data reshape,

together with their syntax, are as follows:

`df.transpose()`: Transposing rows and columns.

`df.stack()`, `df.unstack()`: Stacking and unstacking.

- **Filtering and subset selection:** Some methods for data filtering and subset selection are as follows:

`df.loc[]`, `df.iloc[]`: Selecting subsets.

`df.query()`: Filtering data using a query.

`df.isin()`: Checking for values in a DataFrame.

`df.nlargest()`, `df.nsmallest()`: Selecting the largest or smallest values.

- **Sorting:** Some methods used for sorting are as follows:

`df.sort_values()`: Sorts a DataFrame by one or more columns.

Ascending or descending order.

`df.sort_index()`: Sorts a DataFrame based on the row index.

`sort()`: Sorts lists in ascending and descending order

- **String manipulation:** Some methods used for string manipulation are as follows:

`str.strip()`, `str.lower()`, `str.upper()`, `str.replace()`

Moreover, adding new columns, variables, statistical modeling, testing and probability distribution, and exploratory data analysis is also part of data wrangling and manipulation, which will be covered in [Chapter 2, Exploratory Data Analysis](#).

Conclusion

Statistics provides a structured framework for understanding and interpreting the world around us. It empowers us to gather, organize, analyze, and interpret information, thereby revealing patterns, testing hypotheses, and informing decisions. In this chapter, we

examined the foundations of data and statistics: from the distinction between qualitative (descriptive) and quantitative (numeric) data to the varying levels of measurement—nominal, ordinal, interval, and ratio. We also considered the scope of analysis in terms of the number of variables involved—whether univariate, bivariate, or multivariate—and recognized that data can originate from diverse sources, including surveys, experiments, and observations.

We explored how careful data collection methods—whether sampling from a larger population or studying an entire group—can significantly affect the quality and applicability of our findings. Ensuring data quality is key, as the validity and reliability of statistical results depend on accurate, complete, and consistent information. Data cleaning addresses errors and inconsistencies, while data wrangling and manipulation techniques help us prepare data for meaningful analysis.

By applying these foundational concepts, we establish a platform for more advanced techniques. In the upcoming Chapter 2, Exploratory data analysis we learn to transform and visualize data in ways that reveal underlying structures, guide analytical decisions, and communicate insights effectively, enabling us to extract even greater value from data.

1 Source: https://scikit-learn.org/stable/datasets/toy_dataset.html#iris-dataset

2 Source: <https://github.com/jdorman/awesome-json-datasets#nobel-prize>

3 Source: <https://github.com/jdorman/awesome-json-datasets#nobel-prize>

Join our book's Discord space

Join the book's Discord Workspace for Latest updates,
Offers, Tech happenings around the world, New Release and
Sessions with the Authors:

[**https://discord.bpbonline.com**](https://discord.bpbonline.com)

